

Regularized Dual Averaging

12M42340 チョウ シホウ

7月3日

概要

今回は論文 [1] によって提案された Regularized Dual Averaging (正則化付き双対平均法) を紹介する.

1 導入

1.1 Regularized Stochastic Learning Problem

Regularized stochastic learning (正則化付き確率的学習) 問題とは

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \{ \phi(\mathbf{w}) := \mathbb{E}_z f(\mathbf{w}, z) + \Psi(\mathbf{w}) \} \quad (1.1)$$

ただし, $\mathbf{z} = (\mathbf{x}, y)$ はある未知な分布に従う入力・出力データペア. $f(\mathbf{w}, z)$ は \mathbf{w} と x を利用して y を予測する損失関数 (Loss function) である. なお, 次の幾つの仮定を満たされれば, 最適解を求められる.

- $\Psi(\mathbf{w})$ が閉凸関数.
- $\text{dom } \Psi = \{ \mathbf{w} \in \mathbb{R}^n | \Psi(\mathbf{w}) < +\infty \}$ が閉集合.
- $f(\mathbf{w}, z)$ が凸かつ $\text{dom } \Psi$ における劣微分可能.

損失関数 $f(\mathbf{w}, z)$ の例:

- Least-squares: $x \in \mathbb{R}^n, y \in \mathbb{R}$, and $f(\mathbf{w}, (\mathbf{x}, y)) = \|y - \mathbf{w}^\top \mathbf{x}\|_2^2$
- Hinge loss: $x \in \mathbb{R}^n, y \in \{+1, -1\}$, and $f(\mathbf{w}, (\mathbf{x}, y)) = \max\{0, 1 - y(\mathbf{w}^\top \mathbf{x})\}$
- Logistic regression: $x \in \mathbb{R}^n, y \in \{+1, -1\}$, and $f(\mathbf{w}, (\mathbf{x}, y)) = \log(1 + \exp(-y(\mathbf{w}^\top \mathbf{x})))$

正則化項 (regularization term) $\Psi(\mathbf{w})$ の例:

- ℓ_1 -regularization: $\Psi(\mathbf{w}) = \lambda \|\mathbf{w}\|, \lambda > 1.$
- ℓ_2 -regularization: $\Psi(\mathbf{w}) = \frac{\sigma}{2} \|\mathbf{w}\|_2^2, \sigma > 0.$
- Convex Constraints: $\Psi(\mathbf{w}) = \begin{cases} 0 & \mathbf{w} \in C \\ +\infty & \text{otherwise} \end{cases}$

1.2 一般的なアプローチ

一般的なアプローチは、有限な独立観測値列 z_1, z_2, \dots, z_T を用いて $\phi(\mathbf{w})$ の期待値を近似する。

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{T} \sum_{t=1}^T f(\mathbf{w}, z_t) + \Psi(\mathbf{w}) \quad (1.2)$$

しかし、汎用アルゴリズムを用いて式 (1.2) を直接に解くのは難しい、問題 (1.1) を効率よく解くために Regularized Dual Averaging(RDA) が提案されていた。RDA の本質とは

Regularization + Stochastic Gradient + Dual Averaging + Shrinking

の組み合わせアルゴリズムである。

2 Stochastic Gradient Descent(SGD)

2.1 Online Learning

オンライン学習とは、データを1つずつ読み込んで、それまでの学習結果を更新すること。正則化なし損失関数 $\min f(\mathbf{w}, z)$ に対するのオンライン学習最適化を用いたアルゴリズムは次である。

Algorithm 1 Online Learning Algorithm

initialize: 初期点 \mathbf{w}_1 を求める。

for $t = 1, 2, 3, \dots$, **do**

1. 時刻 t に於ける、入力・出力ペア z_t と \mathbf{w}_t を用いて損失関数 $f(\mathbf{w}_t, z_t)$ を計算。
 2. ステップ1の \mathbf{w}_t を利用して予測を行い、予測が間違えば \mathbf{w}_{t+1} を更新する。
-

間違ったら損失が増加する、オンライン学習の目標は累積損失や Regret などの評価関数を最小化すること。

2.2 Regularized online optimization

正則化付きオンライン最適化は、次の Regret 関数

$$R_t(\mathbf{w}) := \sum_{\tau=1}^t (f_{\tau}(\mathbf{w}_{\tau}) + \Psi(\mathbf{w}_{\tau})) - \sum_{\tau=1}^t (f_{\tau}(\mathbf{w}) + \Psi(\mathbf{w})) \quad (2.1)$$

を使う。

2.3 Stochastic Gradient Descent

次の正則化関数を考え： $\Psi(\mathbf{w}) = I_C(\mathbf{w}) + \psi(\mathbf{w})$ 、 I_C は hard regularization, $\psi(\mathbf{w})$ は soft regularization. SGD は

$$\mathbf{w}_{t+1} = \Pi_C(\mathbf{w}_t - \alpha_t(\mathbf{g}_t + \xi_t)) \quad (2.2)$$

である。ただし, $\mathbf{g}_t \in \partial f(\mathbf{w}_t, z_t), \boldsymbol{\xi}_t \in \partial \Psi(\mathbf{w}_t)$, $\Pi_C(\mathbf{w})$ は次に定義される関数である。

$$\Pi_C(\mathbf{w}) = \arg \min_{\mathbf{w}' \in C} \|\mathbf{w} - \mathbf{w}'\|_2 \quad (2.3)$$

$\Pi_C(\mathbf{w})$ を Euclidean Projection という。

SGD アルゴリズムでは, \mathbf{w}_t と $\phi(\mathbf{w}_t)$ が反復 t 前の情報 $\{z_1, z_2, \dots, z_{t-1}\}$ に依存する。(1.1) の最適解 \mathbf{w}^* の存在を仮定し, $\phi^* = \phi(\mathbf{w}^*)$ とおくと,

$$\lim_{t \rightarrow \infty} \phi(\mathbf{w}_t) = \phi^*$$

を求めるのは SGD のアプローチである。SGD はオンライン学習版の勾配法である。

SGD は Online Learning の一種として, 収束が早い, 大量のメモリを必要としない, 新たなデータが来たときに, 再学習が容易, 機械学習にはよく使われて大規模データにも対応できる, しかし, SGD が正則化付き問題に対して特化していない, もっと効率よく解きたい。

3 Regularized Dual Averaging

3.1 Algorithm

Algorithm 2 Regularized dual averaging(RDA) method

input:

- a strongly convex function $h(\mathbf{w})$ with modulus 1 on $\text{dom } \Psi$, and $\mathbf{w}_0 \in \mathbb{R}^n$, such that

$$\mathbf{w}_0 = \arg \min_{\mathbf{w}} h(\mathbf{w}) \in \arg \min_{\mathbf{w}} \Psi(\mathbf{w}) \quad (3.1)$$

- a pre-determined nonnegative and nondecreasing sequence β_t for $t \geq 1$

initialize: $\mathbf{w}_1 = \mathbf{w}_0, \bar{\mathbf{g}}_0 = 0$

for $t = 1, 2, 3, \dots$, **do**

1. Given the function f_t , compute a subgradient $\mathbf{g}_t \in \partial f_t(\mathbf{w}_t)$
2. Update the average subgradient $\bar{\mathbf{g}}_t$

$$\bar{\mathbf{g}}_t = \frac{t-1}{t} \bar{\mathbf{g}}_{t-1} + \frac{1}{t} \mathbf{g}_t \quad (3.2)$$

3. Compute the next iterate \mathbf{w}_{t+1} :

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}} \left\{ \langle \bar{\mathbf{g}}_t, \mathbf{w} \rangle + \Psi(\mathbf{w}) + \frac{\beta_t}{t} h(\mathbf{w}) \right\} \quad (3.3)$$

RDA の収束率を保証するために次の不等式を満たす補助関数 (auxiliary function) $h(\mathbf{w})$ を定義する,

$$h(\alpha \mathbf{w} + (1 - \alpha) \mathbf{u}) \leq \alpha h(\mathbf{w}) + (1 - \alpha) h(\mathbf{u}) - \frac{\sigma}{2} \alpha (1 - \alpha) \|\mathbf{w} - \mathbf{u}\|^2 \quad (3.4)$$

h は **Strongly Convex Function** と呼ばれ、次の性質が成立

$$h(\mathbf{w}) \geq h(\mathbf{u}) + \langle \mathbf{s}, \mathbf{w} - \mathbf{u} \rangle + \frac{\sigma}{2} \|\mathbf{w} - \mathbf{u}\|^2, \quad \forall \mathbf{s} \in \partial h(\mathbf{u}) \quad (3.5)$$

$\sigma > 0$ を凸性パラメータ (Convexity parameter) という. $\sigma = 0$ の場合は一般的な凸とは同じ.

RDA の各反復 t では

- ステップ1 : 従来の勾配法や劣勾配法と同じ, 劣勾配を算出する.
- ステップ2 : 平均劣勾配を更新する.
- ステップ3 : RDA で構造した簡単関数を計算し, \mathbf{w}_{t+1} を求める.

$\beta_t = \Theta(\sqrt{t})$ の場合, RDA の収束率は

$$\mathbb{E}\phi(\bar{\mathbf{w}}_t) - \phi^* \leq \mathcal{O}\left(\frac{G}{\sqrt{t}}\right) \quad (3.6)$$

ただし, G は \mathbf{g}_t の劣勾配のノルムの上限, $\bar{\mathbf{w}}_t = \frac{1}{t} \sum_{\tau=1}^t \mathbf{w}_t$.

3.2 RDA Methods with General Convex Regularization

$\sigma = 0$ の $\Psi(\mathbf{w})$ に対して, $\mathcal{O}(\frac{1}{\sqrt{t}})$ の収束率を手に入れるため, 次の正数列 $\{\beta_t\}_{t \geq 1}$ を構造する.

$$\beta_t = \gamma\sqrt{t}, \quad t = 1, 2, 3, \dots \quad (3.7)$$

1. Nesterov's dual averaging method

$\Psi(\mathbf{w})$ を凸集合 C に於ける標示関数に定義する.

$$\Psi(\mathbf{w}) = I_C(\mathbf{w}) = \begin{cases} 0 & \mathbf{w} \in C \\ +\infty & \text{otherwise} \end{cases} \quad (3.8)$$

また, $h(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2$ とおくと, 式 (3.3) は次に変形できる

$$\begin{aligned} \mathbf{w}_{t+1} &= \arg \min_{\mathbf{w}} \left\{ \langle \bar{\mathbf{g}}_t, \mathbf{w} \rangle + \Psi(\mathbf{w}) + \frac{\beta_t}{t} h(\mathbf{w}) \right\} = \arg \min_{\mathbf{w}} \left\{ \langle \bar{\mathbf{g}}_t, \mathbf{w} \rangle + I_C(\mathbf{w}) + \frac{\gamma}{2\sqrt{t}} \|\mathbf{w}\|_2^2 \right\} \\ &= \Pi_C \left(-\frac{\sqrt{t}}{\gamma} \bar{\mathbf{g}}_t \right) = \Pi_C \left(-\frac{1}{\gamma\sqrt{t}} \sum_{\tau=1}^t \mathbf{g}_\tau \right) \end{aligned} \quad (3.9)$$

(3.9) 式が Nesterov(2009) の simple dual averaging とは同じである. 近似型 (closed form) がないが, 効率のよい射影関数が存在する.

特に, $C = \{\mathbf{w} \in \mathbb{R}^n \mid \|\mathbf{w}\|_1 \leq \delta\}, \delta > 0$ ならば, 正則化項は hard ℓ_1 -regularization になる.

2. Soft ℓ_1 -regularization

$\Psi(\mathbf{w}) = \lambda \|\mathbf{w}\|_1, \lambda > 0$ と $h(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2$ とおくと, 次の

$$\mathbf{w}_{t+1}^{(i)} = \begin{cases} 0 & \text{if } \|\bar{\mathbf{g}}_t^{(i)}\| \leq \lambda, \\ -\frac{\sqrt{t}}{\gamma} (\bar{\mathbf{g}}_t^{(i)} - \lambda \text{sign}(\bar{\mathbf{g}}_t^{(i)})) & \text{otherwise,} \end{cases} \quad i = 1, 2, \dots, n \quad (3.10)$$

ただし, $\mathbf{w} > 0$ ならば $\text{sign}(\mathbf{w}) = 1$, $\mathbf{w} \leq 0$ ならば $\text{sign}(\mathbf{w}) = 0$,

3.3 RDA Methods with Strongly Convex Regularization

もし $\Psi(\mathbf{w})$ は $\sigma > 0$ の Strongly Convex Function ならば, ある収束率を $\mathcal{O}(\ln t)$ に超えないような非負非減少数列 $\{\beta_t\}_{t \geq 1}$ を構造し, 元問題が $\mathcal{O}(\frac{\ln t}{t})$ のスピードで収束できる. 簡単にするため, $\beta_t = 0$ に設定し, $h(\mathbf{w})$ もいらなくなる.

1. Mixed ℓ_1/ℓ_2^2 -regularization

$\Psi(\mathbf{w}) = \lambda \|\mathbf{w}\|_1 + \frac{\sigma}{2} \|\mathbf{w}\|_2^2, \lambda > 0, \sigma > 0$ とおくと

$$\mathbf{w}_{t+1}^{(i)} = \begin{cases} 0 & \text{if } \|\bar{\mathbf{g}}_t^{(i)}\| \leq \lambda, \\ -\frac{1}{\sigma} \left(\bar{\mathbf{g}}_t^{(i)} - \lambda \text{sign}(\bar{\mathbf{g}}_t^{(i)}) \right) & \text{otherwise} \end{cases} \quad i = 1, 2, \dots, n \quad (3.11)$$

2. Kullback-Leibler(KL) divergence regularization 略

4 Regret bounds and convergence rates

次は式 (2.1) に定義された関数 $R_t(\mathbf{w})$ の上界下界を求める. 便宜上で次の記号を定義する.

$$\Delta_t := (\beta_0 - \beta_1)h(\mathbf{w}_2) + \beta_t D^2 + \frac{L^2}{2} \sum_{\tau=0}^{t-1} \frac{1}{\sigma\tau + \beta_\tau}, \quad t = 1, 2, 3, \dots \quad (4.1)$$

D と L はある定数, σ は $\Psi(\mathbf{w})$ の凸性パラメータ, $\{\beta_\tau\}_{\tau=1}^t$ は RDA の非負非減少入力数列である. また, 収束分析を簡単にするために $\tau = 0$ の時にも値をとれるような $\beta_0 > 0$ となる β_0 を加える. 実は $\beta_0 = \beta_1$ とおけば, (4.1) の $(\beta_0 - \beta_1)h(\mathbf{w}_2)$ を消せる. また, $t = 1$ における \mathbf{w}_2 が分かるため, Δ_1 も求められる. すべての $D > 0$ に対して,

$$\mathcal{F}_D := \left\{ \mathbf{w} \in \text{dom } \Psi \mid h(\mathbf{w}) \leq D^2 \right\} \quad (4.2)$$

に定める.

定理 4.1. *Algorithm 2* で $\{\mathbf{w}_\tau\}_{\tau=1}^t$ と $\{\mathbf{g}_\tau\}_{\tau=1}^t$ を生成し, $\exists L, \forall t \geq 1, \mathbf{g}_t$ の双対ノルム $\|\mathbf{g}_t\|_* \leq L$, が成り立つならば,

$$\forall t \geq 1, \forall \mathbf{w} \in \mathcal{F}_D, R_t(\mathbf{w}) \leq \Delta_t \quad (4.3)$$

もし正則化関数 $\Psi(\mathbf{w})$ の凸性パラメータ $\sigma = 0$ ならば, $\beta_0 = \beta_1$ とおき, (3.7) の $\{\beta_t\}_{t \geq 1}$ を用いて次の Δ_t を導ける.

$$\Delta_t = \gamma \sqrt{t} D^2 + \frac{L^2}{2\gamma} \left(1 + \sum_{\tau=1}^{t-1} \frac{1}{\sqrt{\tau}} \right) \leq \gamma \sqrt{t} D^2 + \frac{L^2}{2\gamma} \left(1 + (2\sqrt{t} - 2) \right) \leq \left(\gamma D^2 + \frac{L^2}{\gamma} \right) \sqrt{t} \quad (4.4)$$

Δ_t を最小化にする γ は $\gamma^* = \frac{L}{D}$, 従って,

$$R_t(\mathbf{w}) \leq 2LD\sqrt{t} \quad (4.5)$$

正則化関数 $\Psi(\mathbf{w})$ の $\sigma > 0$ の場合, $h(\mathbf{w}) = \frac{1}{\sigma} \Psi(\mathbf{w})$ にし, β_t によって違う Δ_t を求められる.

1. Positive constant sequences

$\beta_t = \sigma, t \geq 1, \beta_0 = \beta_1$ に対し

$$\Delta_t = \sigma D^2 + \frac{L^2}{2\sigma} \sum_{\tau=1}^t \frac{1}{\tau} \leq \sigma D^2 + \frac{L^2}{2\sigma} (1 + \ln t) \quad (4.6)$$

2. The logarithmic sequence

$\beta_t = \sigma(1 + \ln t), t \geq 1, \beta_0 = \sigma$ に対し

$$\Delta_t = \sigma(1 + \ln t) D^2 + \frac{L^2}{2\sigma} \left(1 + \sum_{\tau=1}^{t-1} \frac{1}{\tau + 1 + \ln \tau} \right) \leq \left(\sigma D^2 + \frac{L^2}{2\sigma} \right) (1 + \ln t) \quad (4.7)$$

3. The zero sequence

$\beta_t = 0, t \geq 1, \beta_0 = \sigma, h(\mathbf{w}) = \frac{1}{\sigma} \Psi(\mathbf{w})$ に対し

$$\Delta_t \leq \Psi(\mathbf{w}_2) + \frac{L^2}{2\sigma} \left(1 + \sum_{\tau=1}^t \frac{1}{\tau} \right) \leq \frac{L^2}{2\sigma} (6 + \ln t) \quad (4.8)$$

定理 4.2. もし問題 (1.1) には最適解 \mathbf{w}^* が存在し, $\exists D > 0, h(\mathbf{w}^*) \leq D^2$ と $\exists L > 0, \mathbb{E}\|\mathbf{g}\|_*^2 \leq L^2, \forall \mathbf{g} \in \partial f(\mathbf{w}, z), \mathbf{w} \in \text{dom } \Psi$ を満たす, すべて $t \geq 1$ に対して,

$$\mathbf{E}_\phi(\bar{\mathbf{w}}_t) - \phi(\mathbf{w}^*) \leq \frac{\Delta_t}{t}, \text{ where } \bar{\mathbf{w}}_t = \frac{1}{t} \sum_{\tau=1}^t \mathbf{w}_\tau \quad (4.9)$$

参考文献

- [1] L. Xiao, "Dual averaging methods for regularized stochastic learning and online optimization," Journal of Machine Learning Research, vol. 11, pp. 2543-2596, Dec. 2010. [Online]. Available: <http://portal.acm.org/citation.cfm?id=1756006.1953017>
- [2] Yu. Nesterov, Primal-dual subgradient methods for convex problems. Mathematical Programming, 120(1):221-259, 2009.
- [3] S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan, "Stochastic convex optimization," in Proceedings of the Conference on Learning Theory (COLT), 2009. [Online]. Available: <http://eprints.pascal-network.org/archive/00005408/>
- [4] Y. Nesterov, Introductory Lectures on Convex Optimization: A Basic Course (Applied Optimization), 1st ed. Springer Netherlands.
- [5] Steve Wright, NIPS Tutorial, 6 December 2010, <http://pages.cs.wisc.edu/~swright/nips2010/sjw-nips10.pdf>